

## Lec 7

### \* clustering

↳ Problem of finding a hidden structure within unlabeled data.

الفكرة هنا انك بتقسم ال (data) لجزءات، كل جزء هو متشابهة لبعضها بتكون (cluster).

### \* K-means clustering

ex

↳ height, weight and average lifespan of animals.

used for: clustering numerical data.

Input: numerical

↳ Euclidian distance: distance must be ~~existed~~ defined over variable space.

output: Centroid

الباريس مركز كل (cluster)

~~has contents~~

### use cases

→ introduction to classification (Discover classes)

→ exploratory technique

↳ Discover structure in data.

↳ Summarize properties of each cluster.

## Algorithm

① له ادلا حاجة بتختار  $K$  بشكل عشوائي (3 نقط مثلاً)

② بعد كده بتسوي القيم للأقرب لكل نقطة منهم فيعملوا (cluster) مع بعض.

③ هتعمل إعادة حساب للـ (centroid) عشوائية واحدة جديدة.

له هتسب الـ (mean) لكل النقط الموجودة فيه الـ (cluster)

④ هتعيد الخطوة 2، 3 لحد سالتا في الـ (Centroids)

لا تتغير.

أشرح الـ (algorithm) كدهور

موجود في صفحة ١٧ و ٢٨ في معاينة

الـ (output) اللي بتطلع بيها

له الـ (center) بتاع آخر (cluster).

له ~~الـ (center)~~ آخر فرقة بتفرقة له الـ (dataset) في آخر (cluster).

## Picking K

**Heuristic** → Find the "elbow" of within-sum-of squares (WSS), Plot it as function of  $K$

$$WSS = \sum_{i=1}^K \sum_{j=1}^{n_i} |X_{ij} - C_i|^2$$

{

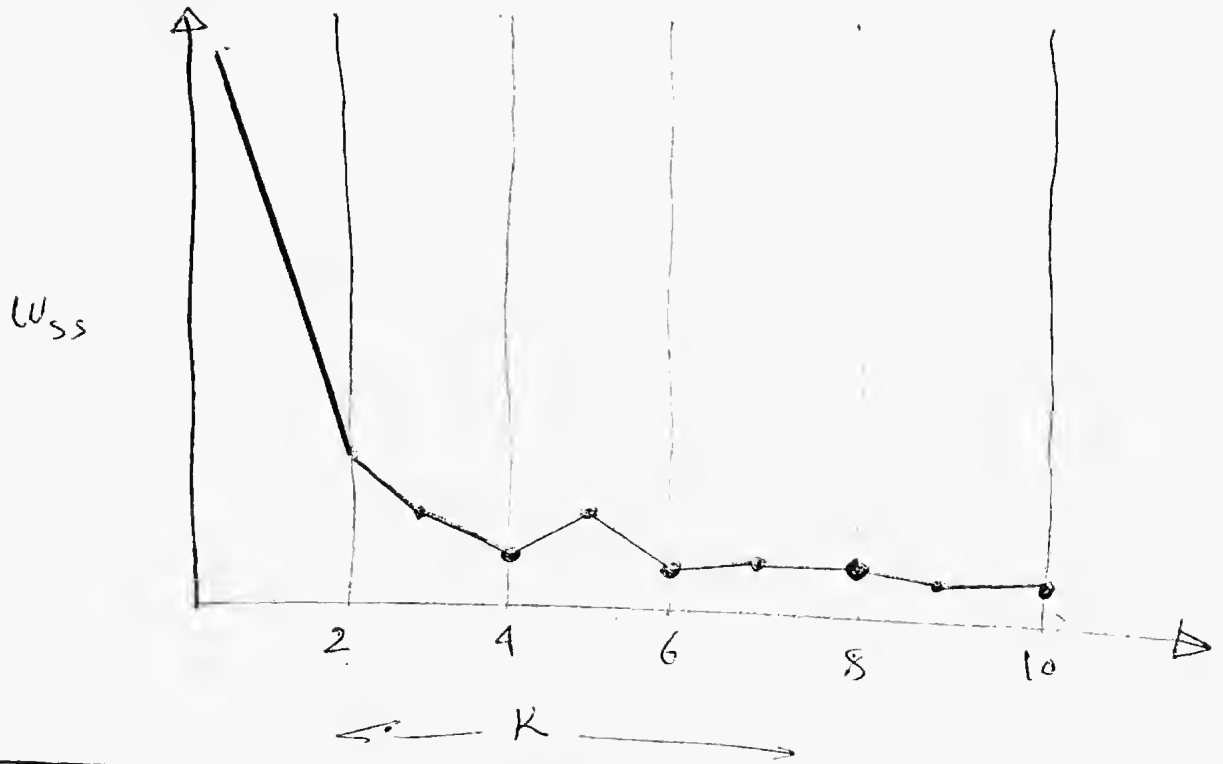
$K \rightarrow$  no. of clusters

$n_i \rightarrow$  no. of points in  $i$ th cluster

$C_i \rightarrow$  centroid of " "

$X_{ij} \rightarrow j$ th point of " "

"Elbows" at  $K = "2, 4, 6"$



مع كل نتيجة نعمل (evaluate)  $N$  (model) فتكون:

(1) هل ال (clusters) باين انها منمهمة عالأقل في بعض الرسومات ولا لا ؟

(2) هل عندك (clusters) فيها (Points) قليلة -  
له جرب تقال قيمة  $K$ .

(3) هل فيه (Centroids) قريبة من بعضها -  
له جرب تقال قيمة  $K$ .

## K-means clustering

Reasons to choose	Cautions.
Easy to implement	Doesn't handle categorical variables.
• Easy to assign new data to existing clusters. ↳ which is the nearest cluster center	↳ sensitive to initialization (first guess)
→ Concise output ↳ coordinates the $K$ cluster centers.	• not scale invariant ↳ variables should all be measured on similar or compatible scales
	↳ Not always desirable ↳ tends to produce "round" equi-sized clusters.

## Lec: 8

### Association rules

↳ unsupervised learning method

↳ used to discover relationships within data.

قواعد الارتباط (Association rules) ←  
التي تكتشف العلاقات في قواعد البيانات (databases)

→ Apriori → works on frequent itemset (set of items that appears together).

### \* Apriori Property

↳ Any subset of a frequent itemset is also frequent.

مثال {shoes, Purses} ← itemset  
وهو يقول أن  $50\% = (\text{Support})$

لأن  $50\%$  من (transactions) لها (itemset)  $50\%$   
لأن يبقى  $50\%$  نقول أن Frequent itemset

→ if  $50\%$  of itemsets have {shoes, Purses} in them then at least  $50\%$  of transactions with have either {shoes} or {Purses} in them → This is Apriori Property

# Lift & Leverage

## Lift

بموجب عدد المرات التي  $x, y$  سيظهروا فيها معاً أكثر من المتوقع لوهما غير معتمدين على بعض.

→ هذه قياس لكيفية  $x, y$  ليها علاقة ببعضهم عن كونهم سيظهروا مع بعض.

$$\text{Lift}(x \rightarrow y) = \frac{\text{support}(x \cap y)}{\text{support}(x) * \text{support}(y)}$$

## Leverage

الفرق بين احتمالية ظهور  $x, y$  معاً و احتمالية

إدراك  $x, y$  غير معتمدين على بعض (وظهروا معاً بشكل متفاجئ)

$$\text{Leverage}(x \rightarrow y) = \text{support}(x \cap y) - \text{support}(x) * \text{support}(y)$$

## Associative rules implementations

1) Market basket analysis

2) recommender systems.

3) discovering web usage patterns.

مثال شرح 11 (associations) rules 17 في مثال 2

رقم 1 في مثال 2 في مثال 3 في مثال 4

1) does data make sense? (data) (check) (data)

2) make "test-set" from hold-out data.

3) Evaluate rules by lift or Leverage.

### Notes

\* support → Percentage of transactions that contain L (set of items)

\* Confidence → Percentage of transactions that contain X, which also contain Y.

\* output of apriori algorithm → set of all rules  $X \rightarrow Y$  with minimum support & confidence.

## \* Apriori

Reasons to choose

→ Easy to implement

→ uses clever observation to Prune Search space (Apriori Property)

→ Easy to Parallelize

Cautions

\* requires many database Scans

\* EXponential time Complexity.

→ Addressed with Lift and Leverage measures

~~Comments~~

8